

Overview of Storage and Indexing

Chapter 6

Data on External Storage

- Disks: Can retrieve random page at fixed cost
 - But reading several consecutive pages is much cheaper than reading them in random order
- File organization: Method of arranging a file of records on external storage.
 - Record id (rid) is sufficient to physically locate record (disk address)
 - Indexes are data structures that allow us to find the record ids of records with given values in index search key fields
- Architecture: Buffer manager stages pages from external storage to main memory buffer pool. File and index layers make calls to the buffer manager.
 - BM is asked fetch a page – diagram on p 20

File Organizations and indexing

A relation is stored as a file of records – records can be inserted, deleted and scanned- step through all the records

- Records are stored in a file in a collection of pages.

Many alternatives exist, *each ideal for some situations, and not so good in others:*

- Heap (random order) files: Suitable when typical access is a file scan retrieving all records.
- Sorted Files: Best if records must be retrieved in some order, or only a 'range' of records is needed.
- Indexes: Data structures to organize records via trees or hashing.
 - Like sorted files, they speed up searches for a subset of records, based on values in certain ("search key") fields
 - Updates are much faster than in sorted files.

Indexes

- An index on a file speeds up selections on the *search key fields* for the index.
 - Any subset of the fields of a relation can be the search key for an index on the relation.
 - *Search key* is *not* the same as *key* (minimal set of fields that uniquely identify a record in a relation).
- An index contains a collection of *data entries*, and supports efficient retrieval of all data entries **k*** with a given key value **k**.
 - Given data entry **k***, we can find record with key **k** in at most one disk I/O. (Details soon ...)

Alternatives for Data Entry k^* in Index

- In a data entry k^* we can store: VERY IMPORTANT SEE page 276
 - Data record with key value k , or
 - $\langle k, \text{rid of data record with search key value } k \rangle$, or
 - $\langle k, \text{list of rids of data records with search key } k \rangle$
- Choice of alternative for data entries is related to the indexing technique used to locate data entries with a given key value k .
 - Examples of indexing techniques: B+ trees, hash-based structures
 - Typically, index contains auxiliary information that directs searches to the desired data entries

Alternatives for Data Entries (Contd.)

- **Alternative 1:** Index is used to store actual data record
 - If this is used, index structure is a file organization for data records (instead of a Heap file or sorted file).
 - Index is used to store actual data record – the order of the index determines the order of the file!
 - At most one index on a given collection of data records can use Alternative 1. (Otherwise, data records are duplicated, leading to redundant storage and potential inconsistency.)
 - If data records are very large, # of pages containing data entries is high. Implies size of auxiliary information in the index is also large, typically.

Alternatives for Data Entries (Contd.)

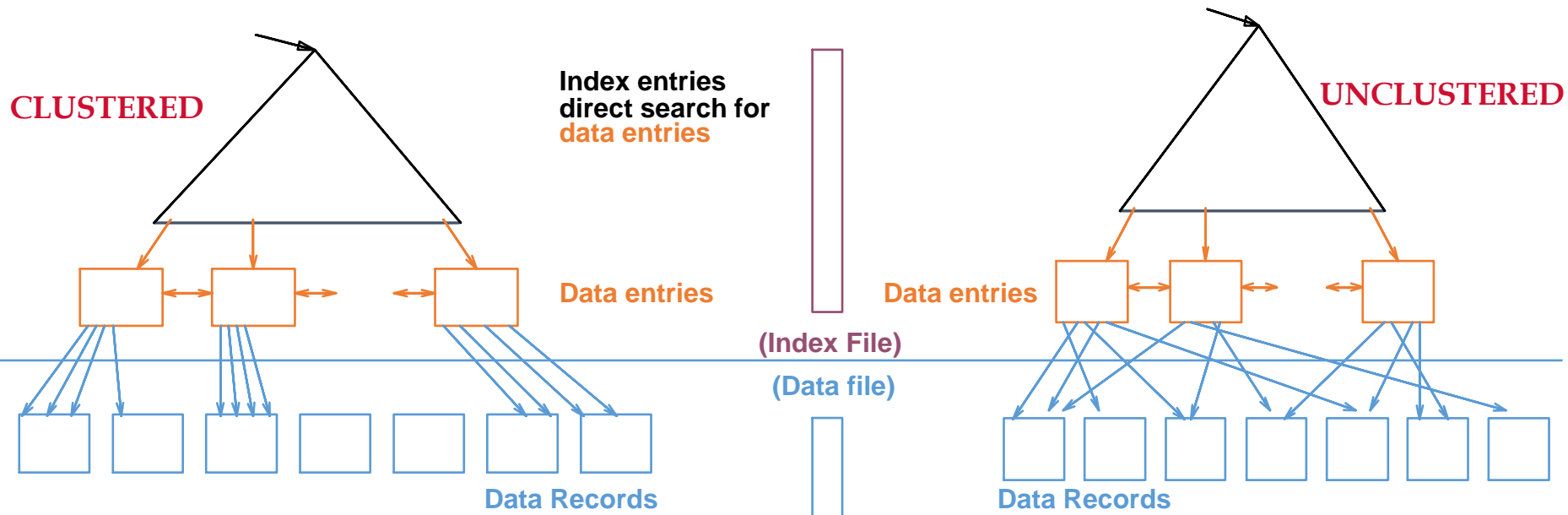
- Alternatives 2 and 3:
- Index is independent of the organization of the file
 - Data entries typically much smaller than data records. So, better than Alternative 1 with large data records, especially if search keys are small. (Portion of index structure used to direct search, which depends on size of data entries, is much smaller than with Alternative 1.)
 - Alternative 3 more compact than Alternative 2, but leads to variable sized data entries even if search keys are of fixed length.
- P277: If there are more than one index on a file – only 1 can be type 1
 - Else we need to store more than 1 copy of the file

Index Classification

- *Primary vs. secondary*: If search key contains primary key, then called primary index.
 - *Unique* index: Search key contains a candidate key.
 - Records are duplicates when they have the same value for the search key
- *Clustered vs. unclustered*: If order of data records is the same as, or 'close to', order of data entries, then called clustered index.
 - Alternative 1 implies clustered; in practice, clustered also implies Alternative 1 (since sorted files are rare).
 - A file can be clustered on at most **one search key**.
 - **The rids in qualifying data entries point to a contiguous collection of records, and we refer to only a few data pages.**
- Cost of retrieving data records through index varies *greatly* based on whether index is clustered or not!

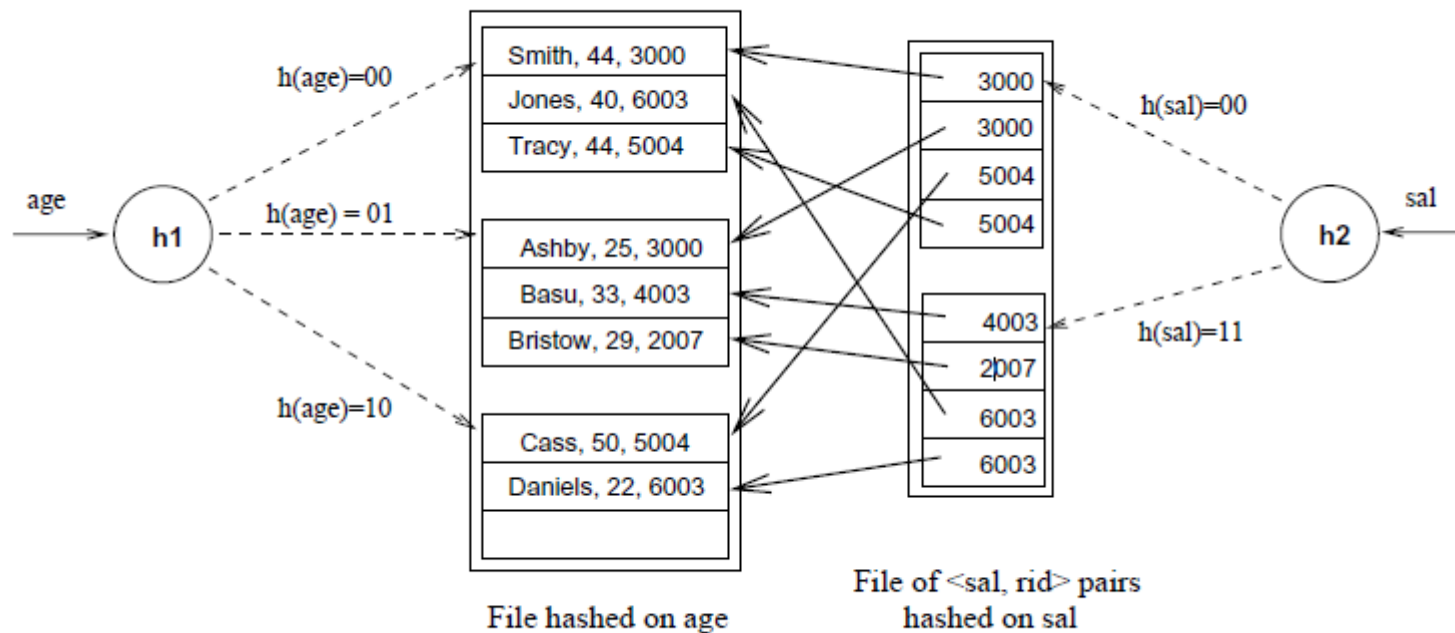
Clustered vs. Unclustered Index NB!!

- Suppose that Alternative (2) is used for data entries, and that the data records are stored in a Heap file.
 - To build clustered index, first sort the Heap file (with some free space on each page for future inserts).
 - Overflow pages may be needed for inserts. (Thus, order of data recs is 'close to', but not identical to, the sort order.)



Hash-Based Indexes (alternative to trees)

- Physical address (page no) is computed using the search key value – so all the records with a certain k^* of “Joe” will be stored in a certain place – bucket.
- Index is a collection of buckets.
 - Bucket = *primary page* plus zero or more *overflow pages*.
 - Buckets contain data entries linked in a chain.
 - Retrieval takes only one or two disk I/Os. (simplify it here to 1 I/O)
- *Hashing function h* : $h(r)$ = bucket in which (data entry for) record r belongs. h looks at the *search key* fields of r .
No need for “index entries” in this scheme.
- Good for equality selections.



Remember, a hashed file is not sorted!

- $h1$ (left) Hashed on age: alternative (1)- the index contains the data itself and file is ordered in the order of the result of the function $h1$. *The function $h1$ converts the age to a binary number and use the 2 least significant bits as rid*
- $h2$ (right) Hashed on salary: alternative (2) – pairs of (key, rid) are stored as index – note equal salaries are in the same page

Tree-based indexing

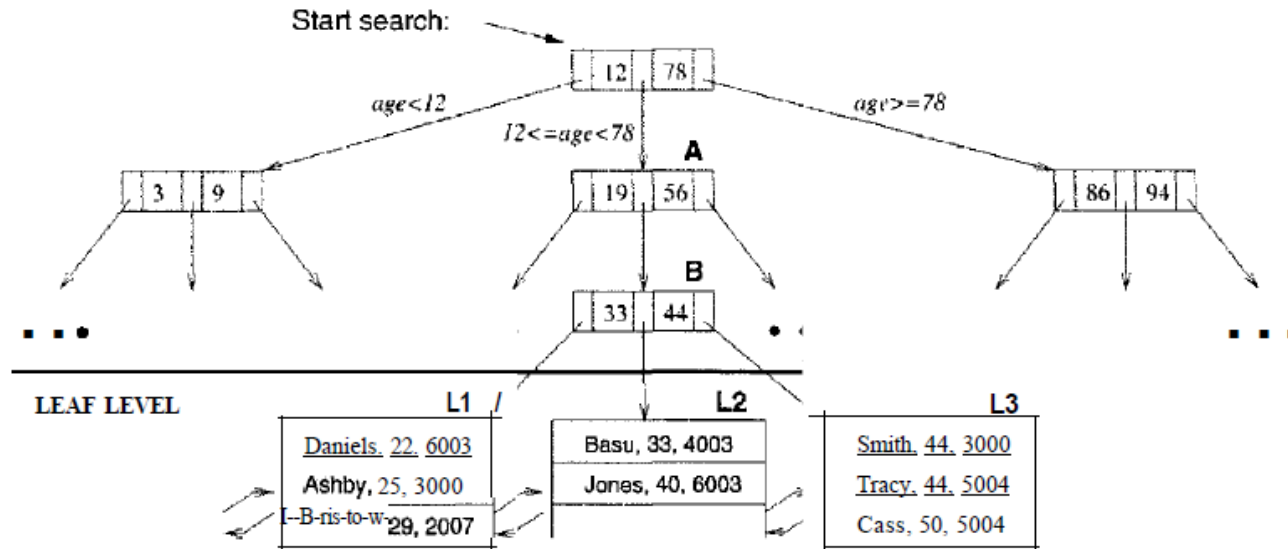


Figure 8.3 Tree-Structured Index

- ❖ Leaf pages contain *data entries*, and are chained (prev & next)
- ❖ Non-leaf pages have *index entries*; only used to direct searches:
 - Diagram shows how to find ages between 24 and 50
 - Leaf pages are linked with double linked list!!!! – So we need only find the first page using the tree!!
 - Number of I/Os is equal to the length of a path from the root to the leaf, plus the number of leaf pages with qualifying data entries

Comparing File Organizations

Assume files are ordered according to composite search key <age,sal>. The following options are tested:

- Heap files (random order; insert at eof)
- Sorted files, sorted on <age, sal>
- Clustered B+ tree file, Alternative (1), search key <age, sal>
- Heap file with unclustered B + tree index on search key <age, sal>
- Heap file with unclustered hash index on search key <age, sal>
- *REMEMBER when a clustered file is searched on another key field, it is similar to searching an unclustered file.*

Operations to Compare

- Scan: Fetch all records from disk
- Equality search
- Range selection
- Insert a record
- Delete a record

Cost Model for Our Analysis

We ignore CPU costs, for simplicity:

- **B:** The number of data pages
- **R:** Number of records per page
- **D:** (Average) time to read or write disk page
 - Measuring number of page I/O's ignores gains of pre-fetching a sequence of pages; thus, even I/O cost is only approximated.
 - Average-case analysis; based on several simplistic assumptions.

C – Comparison of key fields

F- Fan out of tree – avg number of children of a leaf

H – Time to do hashing

** Good enough to show the overall trends!*

Assumptions in Our Analysis

- Heap Files:
 - Equality selection on key; exactly one match.
- Sorted Files:
 - Files compacted after deletions.
- Indexes:
 - Alt (2), (3): data entry size = 10% size of record
 - Hash: No overflow buckets.
 - 80% page occupancy => File size = 1.25 data size
 - Tree: 67% occupancy (this is typical).
 - Implies file size = 1.5 data size

Assumptions (contd.)

- Scans:
 - Leaf levels of a tree-index are chained.
 - Index data-entries plus actual file scanned for unclustered indexes.
- Range searches:
 - We use tree indexes to restrict the set of data records fetched, but ignore hash indexes.

Cost of Operations

| | (a) Scan | (b) Equality | (c) Range | (d) Insert | (e) Delete |
|-------------------------------|----------|--------------|------------|------------|------------|
| (1) Heap | | | | | |
| (2) Sorted | | | | | |
| (3) Clustered | | | | | |
| (4) Unclustered Tree index | | | | | |
| (5) Unclustered Hash index | | | | | |

** Several assumptions underlie these (rough) estimates!*

Comparison formulas:

B - Number of pages to read;

R – Number of records per page;

D – Avg time to read/write disk page - Milliseconds

C – Comparison of key fields – Nanoseconds

F- Fan out of tree – avg number of children of a leaf

H – Time to do hashing - Nanoseconds

- Heap file – unsorted Scan: $B(D+RC)$ Search: $0.5B(D+RC)$
- Sorted file – Search: $D\log_2 B + C\log_2 R$
- Clustered – 67% occupancy: $1.5 * B(D+RC)$

Cost of Operations

| | (a) Scan | (b) Equality | (c) Range | (d) Insert | (e) Delete |
|----------------------------|---------------|-----------------------|--|----------------|----------------|
| (1) Heap | BD | 0.5BD | BD | 2D | Search +D |
| (2) Sorted | BD | $D \log_2 B$ | $D(\log_2 B + \# \text{ pgs with match recs})$ | Search + BD | Search +BD |
| (3) Clustered | 1.5BD | $D \log_F 1.5B$ | $D(\log_F 1.5B + \# \text{ pgs w. match recs})$ | Search + D | Search +D |
| (4) Unclust. Tree index | $BD(R+0.15)$ | $D(1 + \log_F 0.15B)$ | $D(\log_F 0.15B + \# \text{ pgs w. match recs})$ | Search + 2D | Search + 2D |
| (5) Unclust. Hash index | $BD(R+0.125)$ | 2D | BD | Search + 2D | Search + 2D |

** Several assumptions underlie these (rough) estimates!*

Understanding the Workload

- For each query in the workload:
 - Which relations does it access?
 - Which attributes are retrieved?
 - Which attributes are involved in selection/join conditions? How selective are these conditions likely to be?
- For each update in the workload:
 - Which attributes are involved in selection/join conditions? How selective are these conditions likely to be?
 - The type of update (INSERT/DELETE/UPDATE), and the attributes that are affected.

Choice of Indexes

- What indexes should we create?
 - Which relations should have indexes? What field(s) should be the search key? Should we build several indexes?
- For each index, what kind of an index should it be?
 - Clustered? Hash/tree?
- Hash better for equality selection but very poor for range selection – WHY?
- Tree almost similar for equality and range selection...
- Search in tree is faster than search in sorted file
- Never use a clustered file with hashing – WHY?

Choice of Indexes (Contd.)

- **One approach:** Consider the most important queries in turn. Consider the best plan using the current indexes, and see if a better plan is possible with an additional index. If so, create it.
 - Obviously, this implies that we must understand how a DBMS evaluates queries and creates **query evaluation plans!**
 - For now, we discuss simple 1-table queries.
- Before creating an index, must also consider the impact on updates in the workload!
 - **Trade-off:** Indexes can make queries go faster, updates slower. Require disk space, too.

Index Selection Guidelines

- Attributes in WHERE clause are candidates for index keys.
 - Exact match condition suggests hash index.
 - Range query suggests tree index.
 - Clustering is especially useful for range queries; can also help on equality queries if there are many duplicates.
- Multi-attribute search keys should be considered when a WHERE clause contains several conditions.
 - Order of attributes is important for range queries.
 - Such indexes can sometimes enable **index-only** strategies for important queries.
 - For index-only strategies, clustering is not important!
- Try to choose indexes that benefit as many queries as possible.
Since only one index can be clustered per relation, choose it based on important queries that would benefit the most from clustering.

Examples of Clustered Indexes

- B+ tree index on *E.age* can be used to get qualifying tuples.
 - How selective is the condition?
 - Is the index clustered?
- Consider the **GROUP BY** query.
 - If many tuples have *E.age* > 10, using *E.age* index and sorting the retrieved tuples may be costly.
 - Clustered *E.dno* index may be better!
- Equality queries and duplicates:
 - Clustering on *E.hobby* helps!

```
SELECT E.dno
FROM Emp E
WHERE E.age>40
```

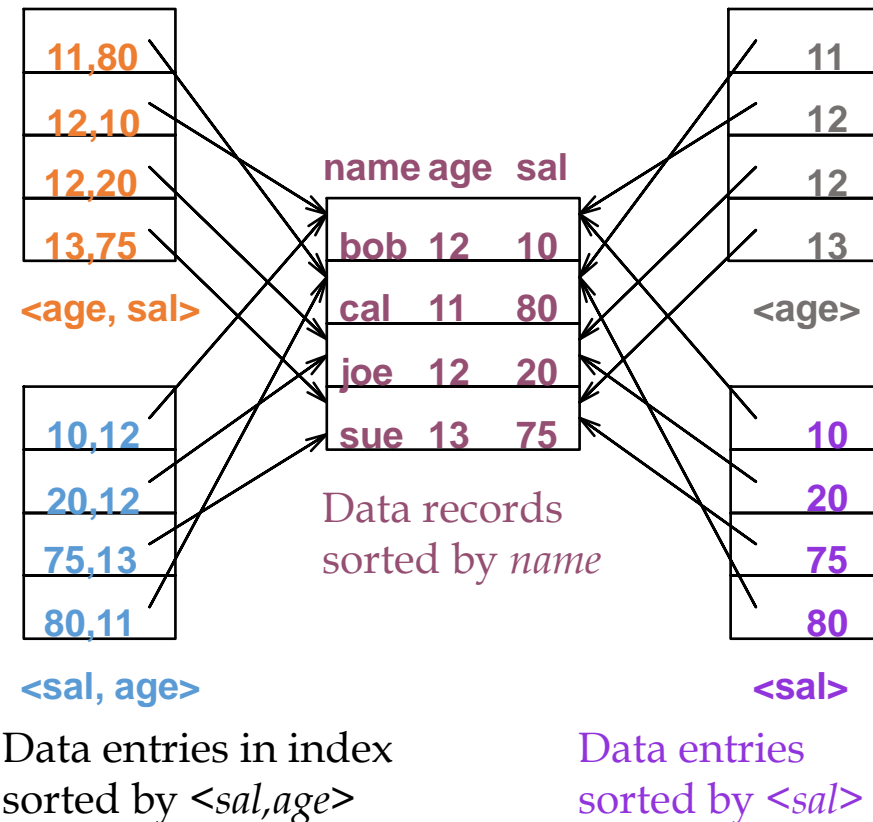
```
SELECT E.dno, COUNT (*)
FROM Emp E
WHERE E.age>10
GROUP BY E.dno
```

```
SELECT E.dno
FROM Emp E
WHERE E.hobby=Stamps
```

Indexes with Composite Search Keys

- **Composite Search Keys:** Search on a combination of fields.
 - **Equality query:** Every field value is equal to a constant value. E.g. wrt $\langle \text{sal}, \text{age} \rangle$ index:
 - age=20 and sal =75
 - **Range query:** Some field value is not a constant. E.g.:
 - age =20; or age=20 and sal > 10
- Data entries in index sorted by search key to support range queries.

Examples of composite key



Composite Search Keys

- To retrieve Emp records with $age=30$ AND $sal=4000$, an index on $\langle age, sal \rangle$ would be better than an index on age or an index on sal .
- If condition is: $20 < age < 30$ AND $3000 < sal < 5000$:
 - Clustered tree index on $\langle age, sal \rangle$ or $\langle sal, age \rangle$ is best.
- If condition is: $age=30$ AND $3000 < sal < 5000$:
 - Clustered $\langle age, sal \rangle$ index much better than $\langle sal, age \rangle$ index!
- Composite indexes are larger, updated more often.
- Ex 8.11 is important!!

Index-Only Plans

- A number of queries can be answered without retrieving any tuples from one or more of the relations involved if a suitable index is available.

<E.dno>

```
SELECT E.dno, COUNT(*)  
FROM Emp E  
GROUP BY E.dno
```

<E.dno,E.sal>

Tree index!

```
SELECT E.dno, MIN(E.sal)  
FROM Emp E  
GROUP BY E.dno
```

<E. age,E.sal>

or

<E.sal, E.age>

Tree index!

```
SELECT AVG(E.sal)  
FROM Emp E  
WHERE E.age=25 AND  
E.sal BETWEEN 3000 AND 5000
```

Index-Only Plans (Contd.)

- Index-only plans are possible if the key is <dno,age> or we have a tree index with key <age,dno>
 - Which is better?
 - What if we consider the second query?

```
SELECT E.dno, COUNT (*)  
FROM Emp E  
WHERE E.age=30  
GROUP BY E.dno
```

```
SELECT E.dno, COUNT (*)  
FROM Emp E  
WHERE E.age>30  
GROUP BY E.dno
```

Index-Only Plans (Contd.)

- Index-only plans can also be found for queries involving more than one table; more on this later.

<E.dno>

```
SELECT D.mgr  
FROM Dept D, Emp E  
WHERE D.dno=E.dno
```

<E.dno,E.eid>

```
SELECT D.mgr, E.eid  
FROM Dept D, Emp E  
WHERE D.dno=E.dno
```

Summary

- Many alternative file organizations exist, each appropriate in some situation.
- If selection queries are frequent, sorting the file or building an *index* is important.
 - Hash-based indexes only good for equality search.
 - Sorted files and tree-based indexes best for range search; also good for equality search. (Files rarely kept sorted in practice; B+ tree index is better.)
- Index is a collection of data entries plus a way to quickly find entries with given key values.

Summary (Contd.)

- Data entries can be actual data records, $\langle \text{key}, \text{rid} \rangle$ pairs, or $\langle \text{key}, \text{rid-list} \rangle$ pairs.
 - Choice orthogonal to *indexing technique* used to locate data entries with a given key value.
- Can have several indexes on a given file of data records, each with a different search key.
- Indexes can be classified as clustered vs. unclustered, primary vs. secondary, and dense vs. sparse. Differences have important consequences for utility/performance.

Summary (Contd.)

- Understanding the nature of the *workload* for the application, and the performance goals, is essential to developing a good design.
 - What are the important queries and updates? What attributes/relations are involved?
- Indexes must be chosen to speed up important queries (and perhaps some updates!).
 - Index maintenance overhead on updates to key fields.
 - Choose indexes that can help many queries, if possible.
 - Build indexes to support index-only strategies.
 - Clustering is an important decision; only one index on a given relation can be clustered!
 - Order of fields in composite index key can be important.